



DO Run II Farms

M. Diesburg, B.Alcorn, J.Bakken, R. Brock, T.Dawson, D.Fagan, J.Fromm,
K.Genser, L.Giacchetti, D.Holmgren, T.Jones, T.Levshina, L.Lueking,
L.Loebel-Carpenter, I.Mandrighenko, C.Moore, S.Naymola, A.Moibenko,
D.Petravick, M.Przybycien, H.Schellman, K.Shepelak, I.Terekhov,
S.Timm, J.Trumbo, S.Veseli, M.Vranicar, R.Wellner, S.White, V.White

DO Farm needs at full rate

- 250K event size
- 50Hz trigger rate
 - peak rate of 12.5 MB/sec
 - DC is less but reprocessing will bring back up

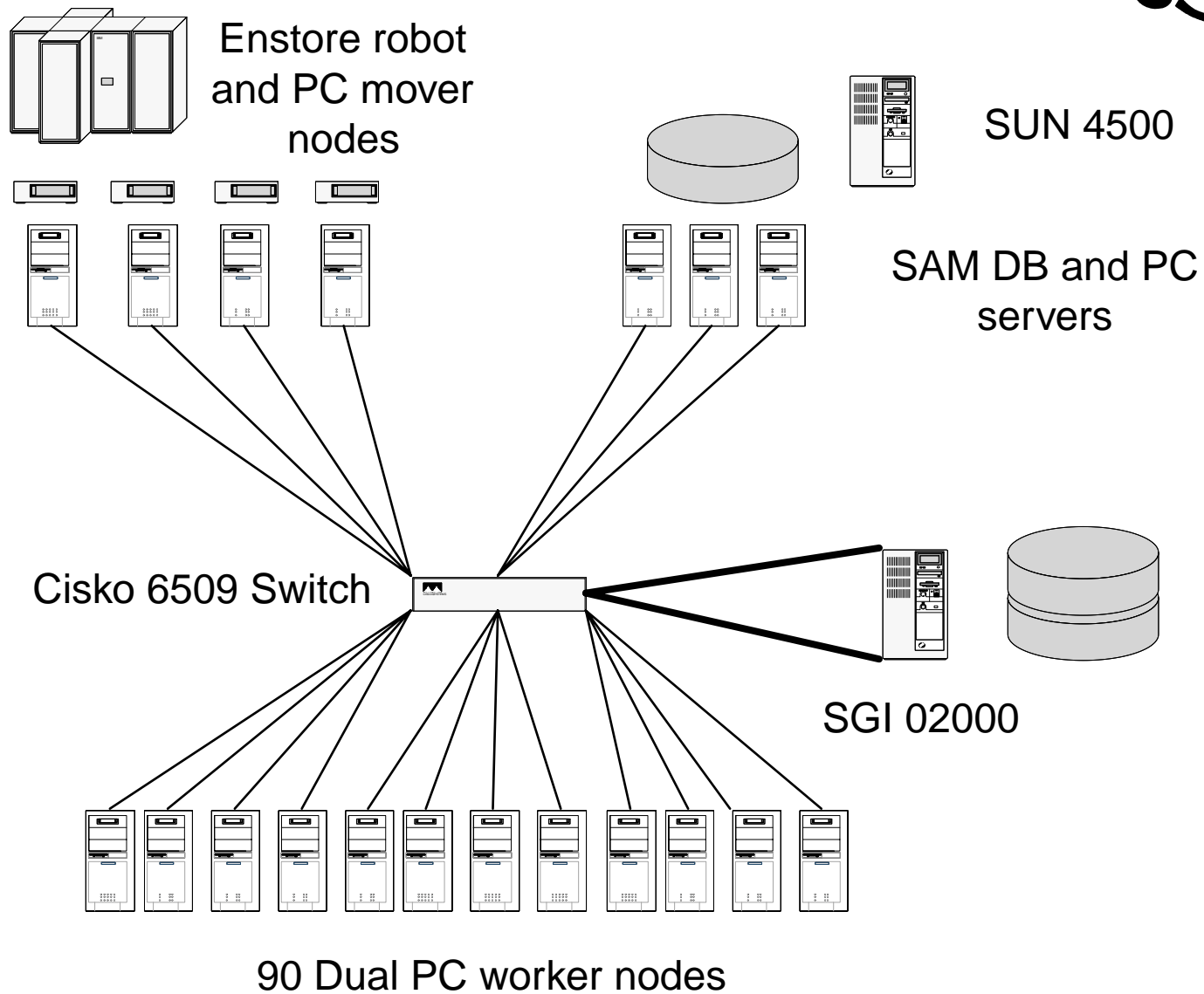
Goal

Reality

- Reconstruction 5- 13 seconds/event

on 750 MHz PIII

- need 250->>500 CPU's to handle peak rate
- DC is 40% of peak
- time constant for 1 GB file is > 12 hours.



I/O machine

- Purpose
 - split/merge of farm output
 - Serve home areas
 - Batch system control
 - File delivery master
- DObbin
 - 4 CPU SGI O2000
 - 2 GB ethernet cards
 - 4 72 GB disk partitions (2 way stripe) + 1 TB raid
 - peak I/O rates of 40-60 MB/sec

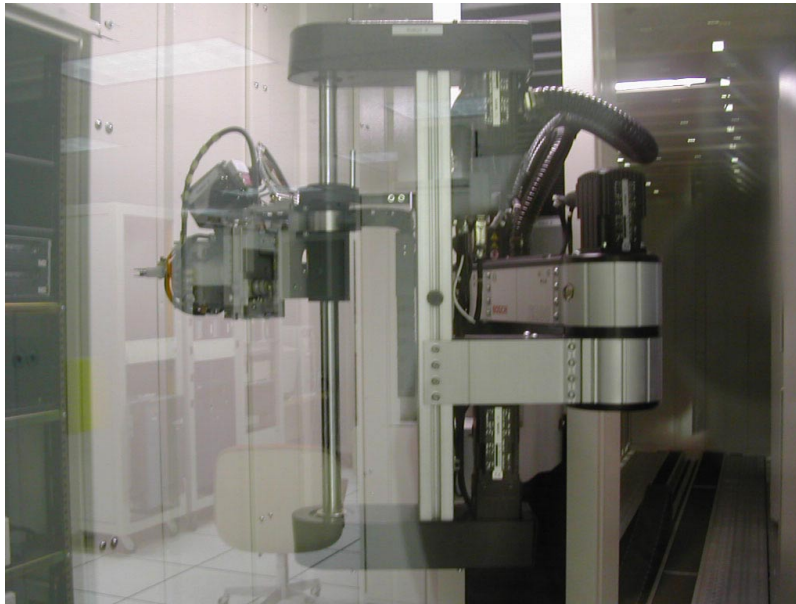


Worker Nodes

- 40 Dual Pentium III 500MHz
 - 256MB/CPU
 - 10 on loan to gtr
 - 10 for test
- 50 Dual Pentium III 750MHz
 - 512 MB/CPU
 - 12 on loan to L3
- 2 data disks (18 GB) + 6GB system
- 100Mb ethernet
- CD/floppy for system configuration



Tape robot



The robot

- Holds 10,000 8mm MII tapes -> 500 TB
- ~ 10 MII tape drives
- 12.5 MB/sec
- Major problems with tape/drive reliability
- Replacement technology on the way

How it works

- Shift captains produce 'shiftsets'
- Farm shifter submits job for shiftset
 - Spread the job across N nodes where $N \sim N_{\text{files}}$
- Data are copied from tape to nodes
- D0 executable tarfile is copied to nodes
- Runs for 12-24 hours
- Data are copied back to I/O node
- Reco files -> tape
- Root files -> merge process -> tape and disk

Tentacles!

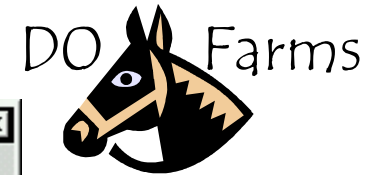


Shared libraries

*Special config files, with
hard-wired paths*

*Databases with expensive
licenses*

Farm Batch System Monitor

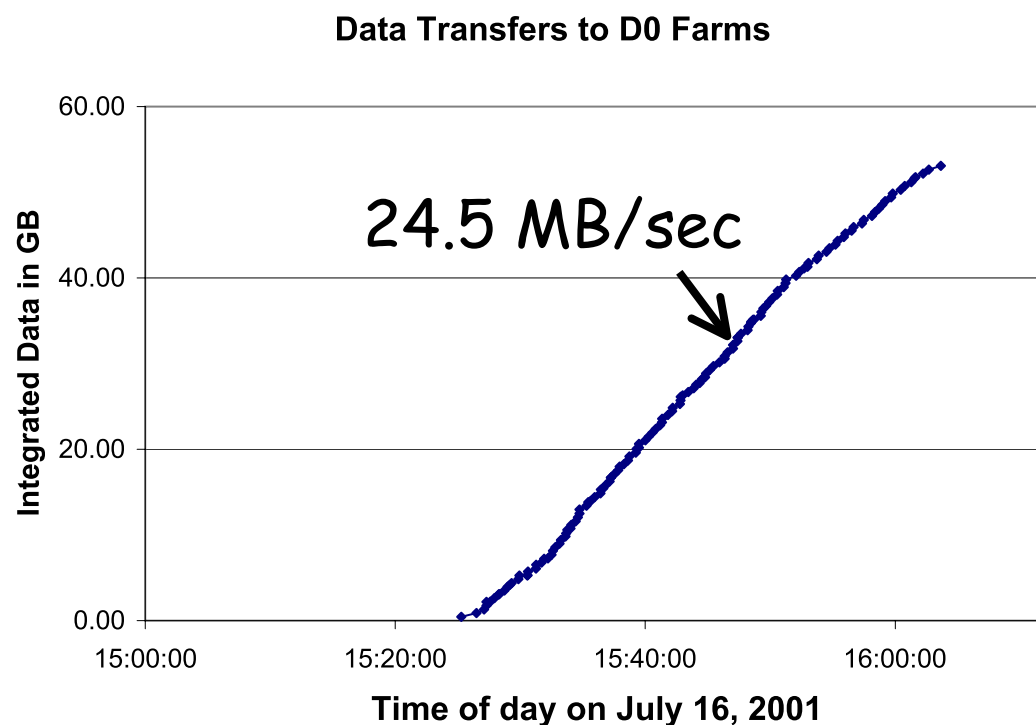


100% of dual

50% of dual CPU

Efficient use of of CPU

I/O Results of typical farm startup



- Cold start of $\frac{1}{2}$ of the D0 farm.
- 90 receiver nodes
- 141 files of average size 376 MB
- Read from 2-3 network mounted Mammoth II tapes over 100 MB ethernet at ~10MB/sec/drive.
- Elapsed time of 44 minutes.
- This is twice peak rate from the detector.

Farm rate calculator

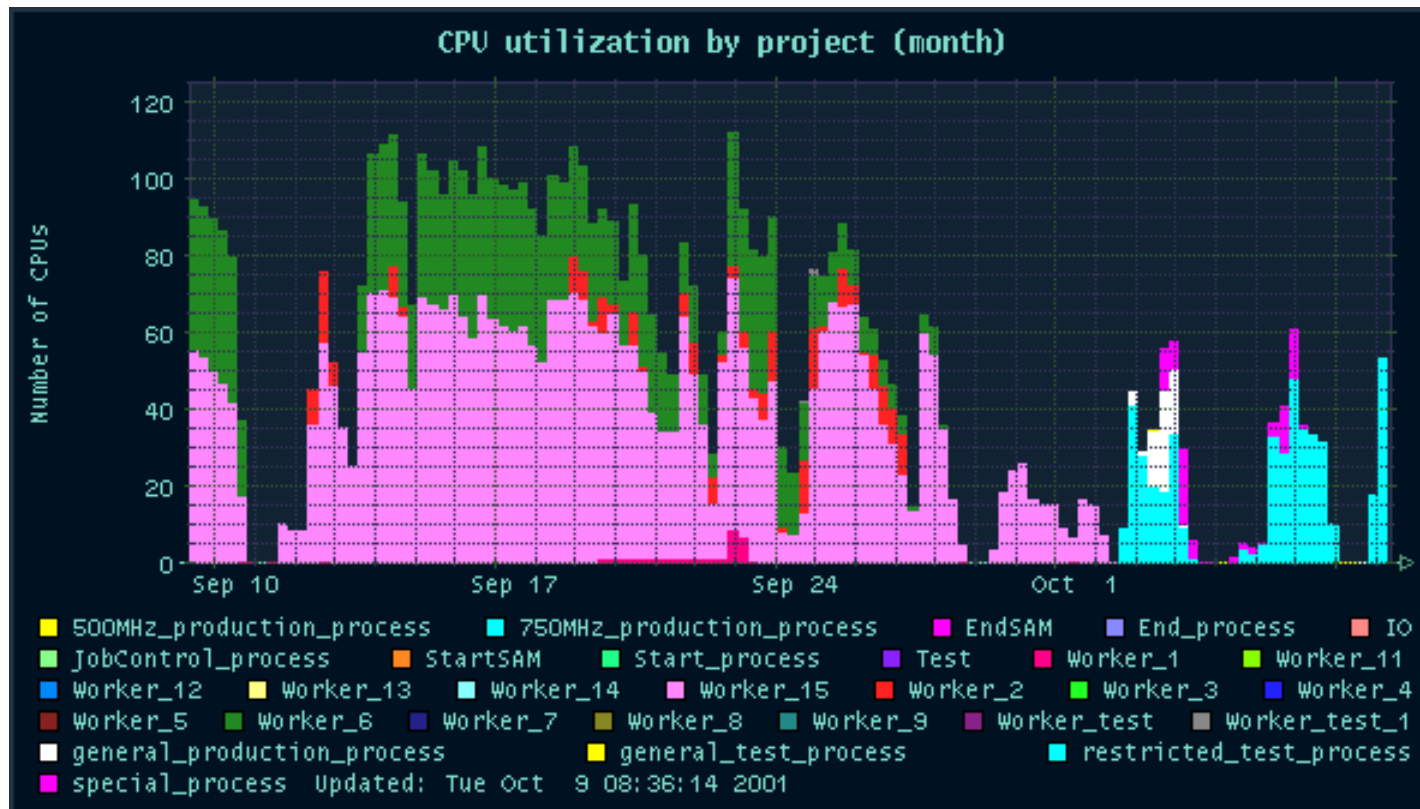
Reco measured to take 13 seconds on 'fast' nodes					recoA is 4 seconds!			
32 new machines with specint95 ~ 40 will arrive soon							efficiency	
		cpu time,seconds				number	1.00	0.70
	specint95	reco	analyze	total		of CPU's	rate, Hz	rate, Hz
Slow	21.0	22.3	6.0	28.3		40	1.4	1.0
Fast	36.0	13.0	3.7	16.7		100	6.0	4.2
new	48.0	9.8	2.8	12.5		64	5.1	3.6
fy2002	60.0	7.8	2.2	10.0		200	20.0	14.0
					DC rate with existing farm			5.2
					DC rate with nodes now on the floor			8.8
					and purchase 100 in fy2000			22.7
					and purchase 100 more			36.7

Rates should be
compared to online
rate \times acc/logging
efficiency



We designed for 50 Hz at
40% efficiency = 20 Hz

Current Production



Overall CPU utilization was 40% for September

We were able to keep up with data taking to within 2-3 days.

Status of recent production

T01.56.00

- Data from 8/25 to 9/24
- Runs 129194-132100
- 7.9 M non-daq_test events
- 6.2 M in shift-sets
- 3.6 M reconstructed
- 3.4 M reco available
- 3.5 M root available

P10.04.00

- Data from 9/25 to 10/7
- Runs 132100 on
- 2.2M non-daq_test events
- 1.8 M in shiftsets
- 1.2 M reconstructed (still running)
- 0.95 M reco available
- 0.9 M root available

Tape Statistics

When tapes are known bad, they are labeled
NOACCESS or NOTALLOWED

At any given time, other tapes may be unreadable

- 10/93 raw tapes are bad, 358/4534 GB.
- 13/118 reco tapes are bad, 586/6390 GB.
- 7/13 root tapes are 48/62GB
- For p10.04.00 the root losses were 0 because of disk backup and merging

Future

- Reprocessing

- Do the recent data with p10.07.00
- 8M events at 17 seconds/event / (125 CPU's * .7 eff)>
- ~ 19 days.

Longer term

- **Upgrades**

- Currently could purchase 1.1 GHz machines, only 40% faster than existing nodes - need 200 with existing code speed -> 300K\$
- 1.8GHz P4 will be available ~January
- Partial purchase followed by big purchase?
- New tape technology on the way.

Status

- System has been in use for MC processing since before CHEP 2000
- System has been processing data as it comes off the D0 detector since March 2001
- Hardware/control/monitoring can handle full data rates well but...
- Major problem is speed of executable and data expansion during detector debugging
 - Output size is \sim input size by design
 - Currently factor of 2-3 larger due to debugging info.
 - Better thresholds and less noise will make life much easier
- *Farms get more stress at beginning of run than later!!*

Structure of a Farm Job

- **START:**
 - Tell SAM which files you will want
 - Go into wait state until get end signal
- **WORKER: runs on N nodes**
 - Download D0 environment
 - Inform SAM ready for data
 - Ask for SAM for next file
 - Process file and store output to output buffer
 - Inform SAM of success and ask for next file
 - On error or end of list, terminate.
- **END:**
 - Create job summary
 - Send message to Start process telling it to shut down the SAM connection for input